

PATENT  
Atty. Dkt. No. YOR920010320US1

**AMENDMENTS TO THE CLAIMS:**

This listing of claims will replace all prior versions, and listings, of claims in the application:

RECEIVED  
CENTRAL FAX CENTER

JUN 18 2008

**LISTING OF CLAIMS**

1. (Currently Amended) A method, in a network comprising a primary server and a plurality of offload servers, for dynamic offloading of processing requests from said primary server to any one of said plurality of offload servers, the method comprising the steps of:

determining a load on said primary server;

if the load on said primary server is less than a first threshold, serving processing requests at said primary server; and

only if the load on said primary server exceeds said first threshold, then offloading at least a portion of said processing requests to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request requests and the offloaded at least a portion of said processing requests is the only work handled by said plurality of offload servers; and

if the load on said primary server exceeds a second threshold, throttling at least one of said processing requests.

2. (Original) The method of claim 1 wherein said load comprises bandwidth utilization and said first threshold is a network bandwidth utilization of said primary server.

3. (Currently Amended) The method of claim 1 wherein the said load comprises CPU utilization and said first threshold is a central processing unit (CPU) utilization of said primary server.

4. (Currently Amended) The method of claim 1 wherein serving the processing

requests at said primary server includes returning a page to a user wherein all the embedded objects in the page have links to said primary server; and

wherein offloading at least a portion of the processing requests to any one of said plurality of offload servers includes serving a base page at said primary server in which the links for embedded objects point to any one of said plurality of offload servers.

5. (Previously Presented) The method of claim 1 wherein offloading at least a portion of the processing requests to any one of said plurality of offload servers includes routing an incoming Web request to a selected offload server.

6. (Cancelled)

7. (Currently Amended) The method of claim [[6]] 1 wherein throttling at least one of said processing request requests includes returning a page to a user indicating that a server is overloaded.

8. (Currently Amended) The method of claim [[6]] 1 wherein throttling at least one of said processing requests includes dropping the at least one of said processing request requests without returning any information to a user.

9. (Currently Amended) The method of claim [[6]] 1 wherein throttling at least one of said processing request requests includes returning a page to a user indicating that a server is overloaded if said load exceeds said second threshold, and dropping said at least one of said processing request requests if said load exceeds a third threshold.

10. (Currently Amended) The method of claim 1 wherein the a determination of which of said plurality of offload servers that at least one a portion of said processing request requests is to be offloaded to is based on one or more of a group including: a client identity, a client gateway (IP) (Internet Protocol) address, a price of the offload service, or a current or previous load on the ~~at least one~~ any one of said plurality of

offload server servers

11. – 31. (Cancelled)

32. (Currently Amended) A method for allocating processing requirements on an IP Internet Protocol network between a primary server and a plurality of offload servers, comprising:

periodically evaluating processing requests to determine a load on said primary server;

if said load exceeds a first threshold, for a predetermined period of time directing at least one processing request to any one of said plurality of offload servers, wherein all of said plurality of offload servers are configured to process said processing request and the ~~offloaded~~ at least a portion of said processing requests is the only work handled by said plurality of offload servers; and

only if said ~~processing~~ load does not exceed said first threshold, directing said processing requests to said primary server; and

if the load on said primary server exceeds a second threshold, throttling at least one of said processing requests.

33. (Currently Amended) The method of claim 32 wherein said load comprises network bandwidth and said first threshold is a measure of the network bandwidth utilization of said primary server.

34. (Currently Amended) The method of claim 32 wherein said load comprises central processing unit (CPU) utilization and said first threshold is a measure of the CPU utilization of said primary server.

35. (Currently Amended) The method of claim 32 wherein directing said processing requests to said primary server further includes returning a page to a user wherein all the embedded objects in the page have links to said primary server; and

PATENT  
Atty. Dkt. No. YOR920010320US1

directing at least one processing request to any one of said plurality of offload servers further includes serving a base page at said primary server in which the links for embedded objects point to said any one of said plurality of offload servers.

36. (Previously Presented) The method of claim 32 wherein directing at least one processing request to any one of said plurality of offload servers further includes routing an incoming Web request to a selected offload server.

37. (Currently Amended) The method of claim 32 ~~and further including, if said load exceeds a second threshold,~~ wherein said throttling at least one of said processing request requests comprises by returning a page to a user indicating that a server is overloaded.

38. (Currently Amended) The method of claim 32 ~~and further including, if said load exceeds a second threshold,~~ wherein said throttling of at least one of said processing requests comprises dropping the at least one of said processing request requests without returning any information to a user.

39. (Currently Amended) The method of claim 32 ~~and further including~~ wherein the throttling of at least one of said processing request requests comprises by returning a page to a user indicating that the primary server is overloaded if the ~~primary server~~ load exceeds ~~[[a]]~~ the second threshold, and further comprising dropping the at least one of said processing request requests if the ~~primary server~~ load exceeds a third threshold.

40. (Currently Amended) The method of claim 32 further including determining which of said plurality of offload servers that said at least one of said processing request requests is to be offloaded to is based on one or more of a group including: a client identity, a client gateway ~~(IP)~~ (Internet Protocol) address, a price of the offload service, or a current or previous load on the ~~at least one~~ any one of said plurality of offload server servers.

PATENT

Atty. Dkt. No. YOR920010320US1

41. - 42. (Cancelled)